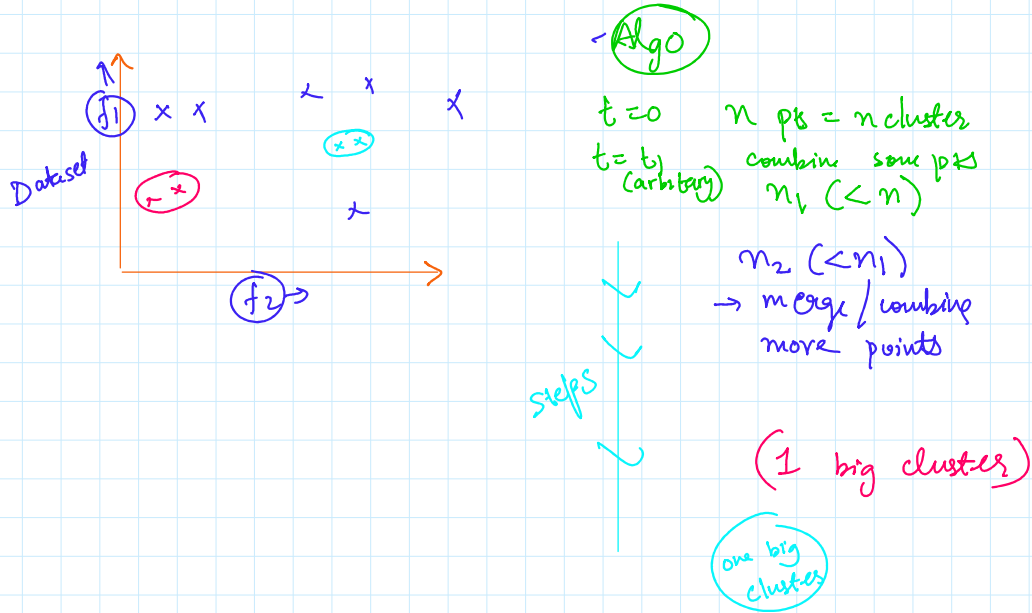
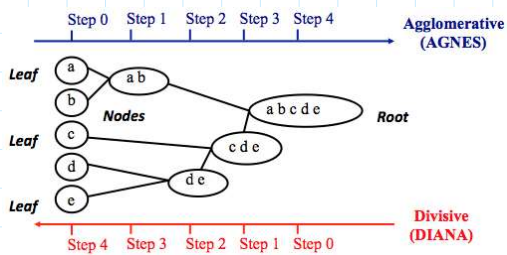
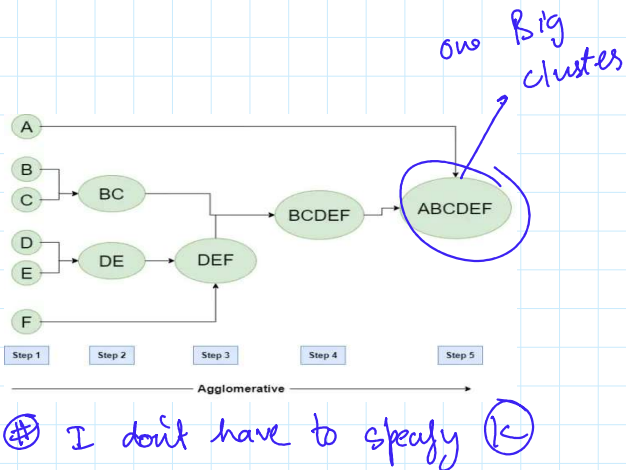
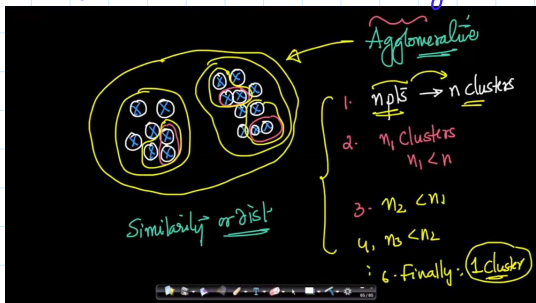


What is Agglomerative Clustering?

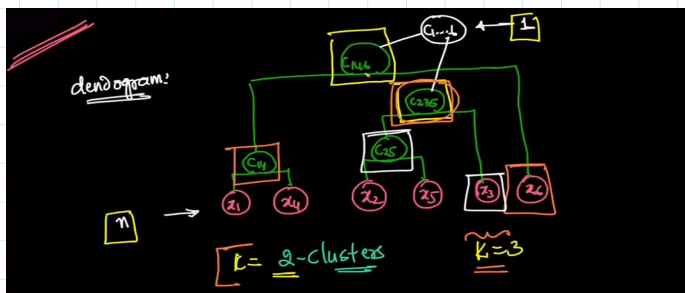
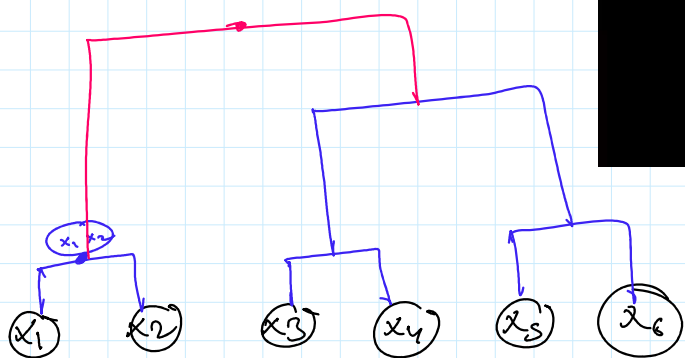
- The word agglomerative suggests combining things
- It is a **bottom-up** approach
- Agglomerative clustering starts with the assumption that every datapoint is a cluster
- Then, it groups the clusters which are closed to each other until there is only a single cluster left



### Agglomerative Clustering



centroid



① distance every point  $(x_i) \rightarrow x_2, x_3, \dots, x_6$   
 $x_2 \rightarrow$  every other pts  
 $** (N C_2)$

Divisive Clustering  $\rightarrow$

- What is Divisive Clustering?
- It is a complete opposite of agglomerative approach
  - It is a **top-down** approach
  - It starts with one big clusters that contains all the datapoints.
  - It then divides the points into different clusters till each data point is a cluster itself

Proximity Matrix  $\rightarrow$  matrix of distance (similarity)

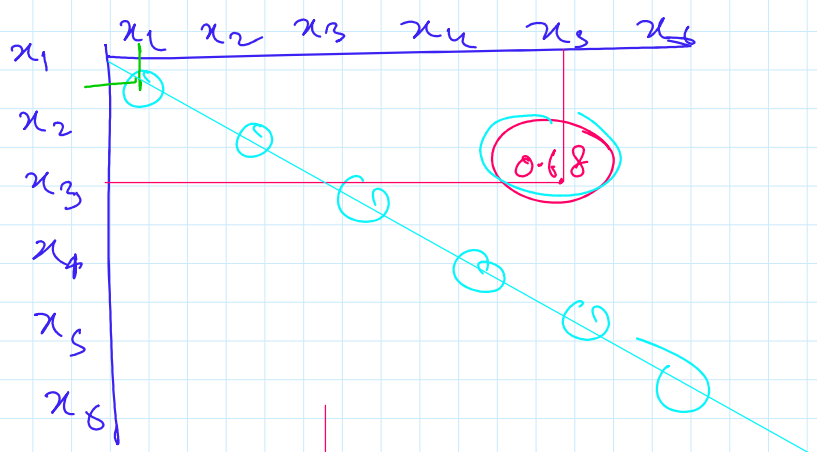
Input = A list of points =  $[x_1, x_2, x_3, x_4, x_5]$

Output =  $n \times n$  Matrix where each  $(i, j)$  represents distance b/w  $i$ th and  $j$ th clusters (or data point)

Small example

Data =  $x_1, x_2, x_3, x_4, x_5, x_6$

(Pairwise Distance)



$$\text{Dist}(x_1, x_2) = \text{Dist}(x_2, x_1)$$

$$* (\text{Dist}(x_3, x_5) = 0.68)$$

merge  $x_3$  and  $x_5$

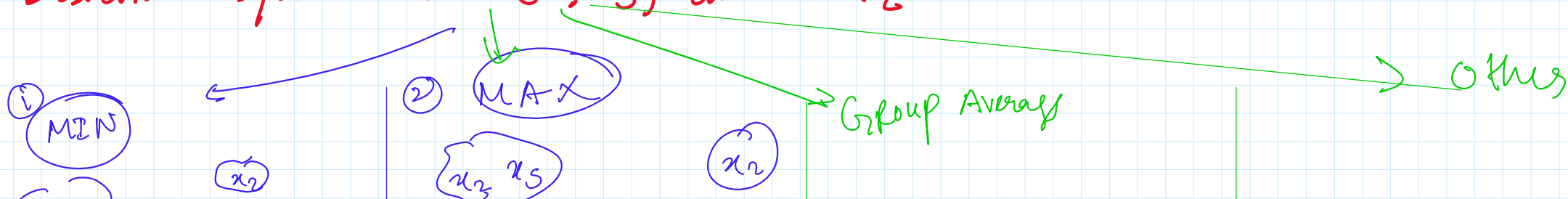
MIN

$C = \{x_3, x_5\}$

$x_1$	$x_2$	$(x_3, x_5)$	$x_4$	$x_6$
0	0.01			
$x_2$				
$x_3$	0.35	0.6	0	0.2
$x_4$				0

$$\text{euc} = x_2, x_1 = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

Distance b/w cluster  $(x_3, x_5)$  and  $x_2$



$$\min \left[ d(x_3, x_2), d(x_5, x_2) \right]$$

$$\max \left[ d(x_3, x_2), d(x_5, x_2) \right]$$

$$\frac{d(x_3, x_2) + d(x_5, x_2)}{2}$$

	$(x_1, x_2)$	$(x_3, x_5)$	$x_4$
$(x_1, x_2)$			
$(x_3, x_5)$			
$x_4$			

## Distance Metrics

Q How to compute  $\text{dist}(C_i, C_j)$ ?

$\rightarrow d$  b/w centroids

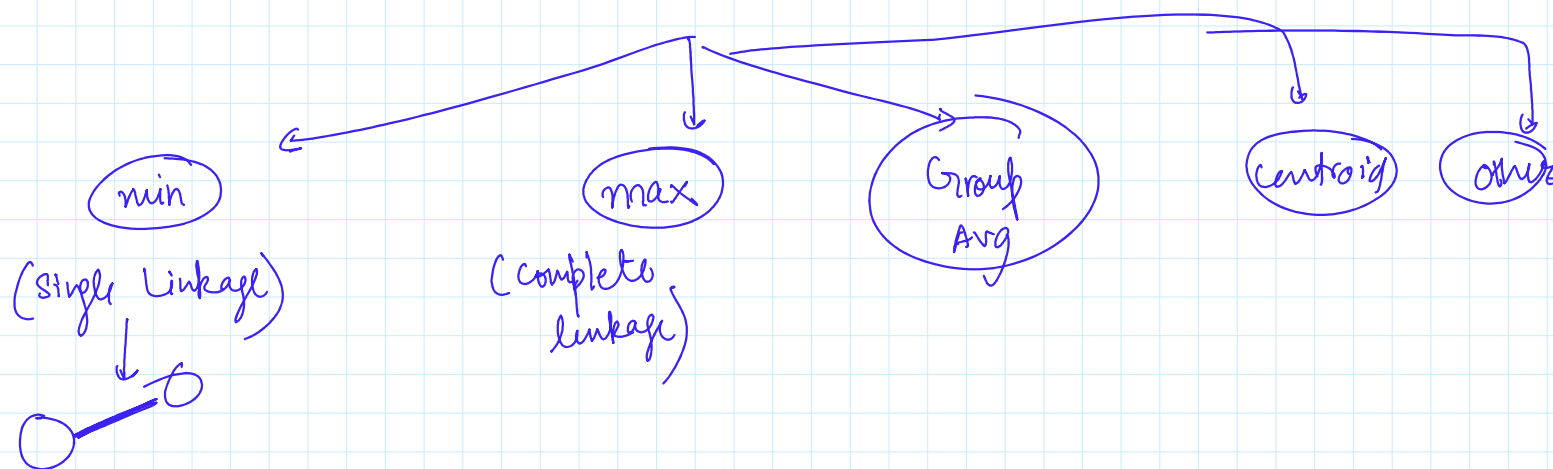
1. euc. dist b/w centroids

2. MAX:  $\max_{x_i \in C_i, x_j \in C_j} \text{dist}(x_i, x_j)$

3. MIN:  $\min_{x_i \in C_i, x_j \in C_j} \text{dist}(x_i, x_j)$

4. group average dist:  $\frac{\sum_{x_i \in C_i} \sum_{x_j \in C_j} \text{dist}(x_i, x_j)}{|C_i| |C_j|}$

5. Ward's dist:  $\frac{\sum_{x_i \in C_i} \sum_{x_j \in C_j} \text{dist}(x_i, x_j)^2}{|C_i| |C_j|}$



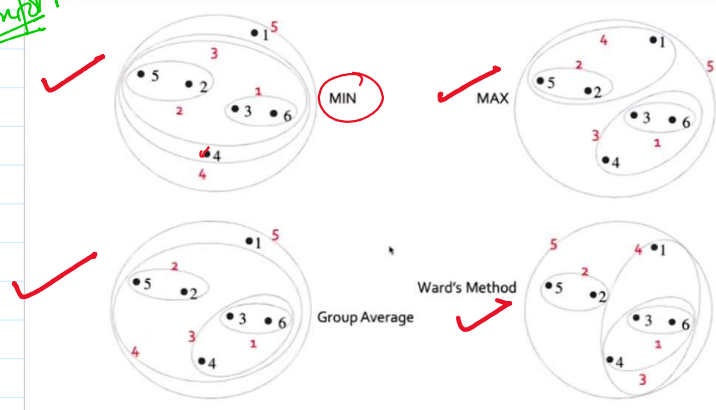
Break  $\rightarrow$  10:18 pm  $\rightarrow$  6 min break

Q How will the distance affect?

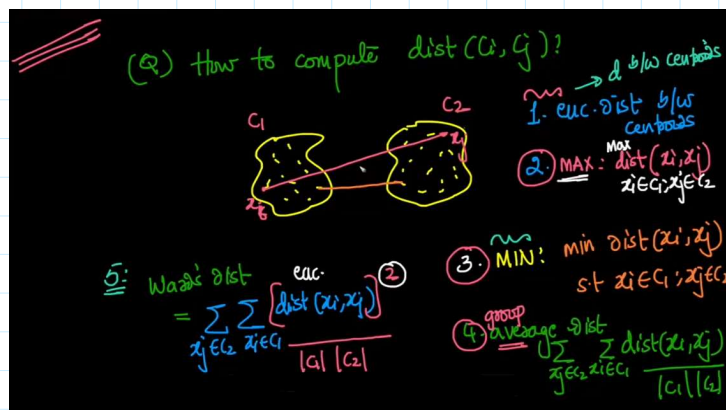
Real  $\rightarrow$  (It depends on the data // Domain specific)

...

No Informat



Ward Distance



$$W.D \rightarrow \frac{\sum_{x_i \in C_1} \sum_{x_j \in C_2} [dist(x_i, x_j)]^2}{|C_1| |C_2|}$$

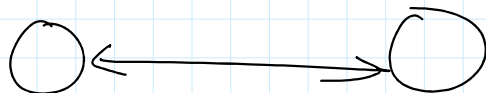
eucledian

$$eucledian^2 = \left( \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2} \right)^2 = d$$

ward Distance  $\Rightarrow$

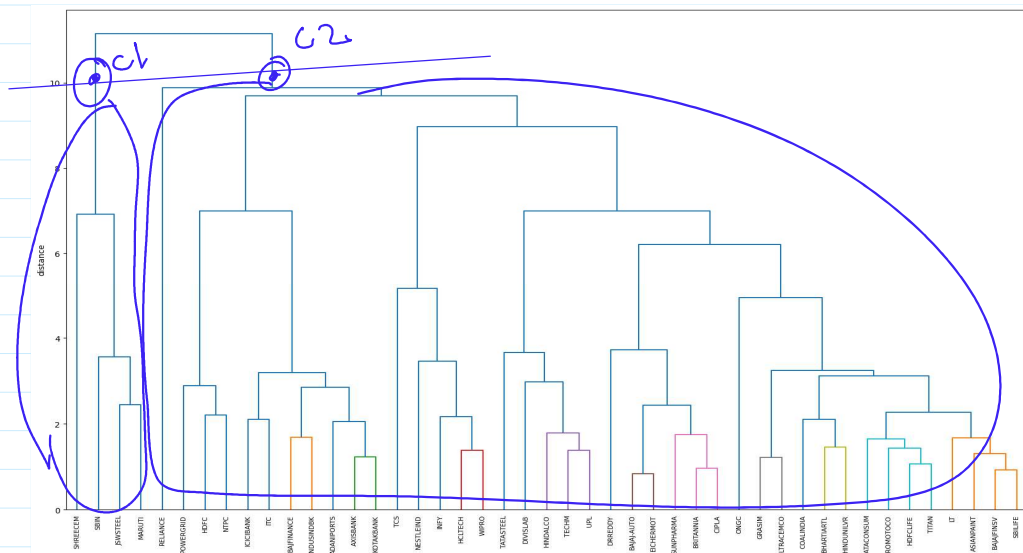
$$\frac{(x_2 - x_1)^2 + (y_2 - y_1)^2}{|C_1| |C_2|}$$

Dissimilarity



Financial Terms

1. Market Cap (Market Capitalization): This is the total market value of a company's outstanding shares of stock. It is calculated by multiplying the current market price of one share by the total number of outstanding shares. It's an indication of the company's size and has implications for how volatile the stock price might be and how broad the ownership might be.
2. Regular Market Volume: This refers to the number of shares of a stock that are traded during regular trading hours on a given day. It is an indicator of the stock's liquidity and the level of investor interest in the company.
3. Quarterly Earnings Growth: This metric shows the percentage increase or decrease in a company's earnings per share (EPS) from one quarter to the next. It is an important indicator of a company's profitability and its potential to increase revenue and manage costs over time.
4. Book Value: The book value of a company is calculated from the balance sheet, and it's the difference between a company's total assets and total liabilities. It represents the net asset value of the company according to accounting principles.
5. Total Revenue: This is the total amount of income generated by the sale of goods or services related to the company's primary operations. It is often referred to as the "top line" because it is the first number you see on a company's income statement.
6. Return on Assets (ROA): This is a profitability ratio that measures how effectively a company can earn a return on its investment in assets. It is calculated by dividing the company's net income by its total assets.
7. Profit Margins: This financial metric measures the percentage of revenue that exceeds the costs associated with a company's products or services. There are various types of profit margins, including gross, operating, and net profit margins, each providing different levels of cost and profit insight.
8. Earnings Growth: This represents the annual compounded growth rate of earnings from investments over time. It indicates the long-term earning potential and financial stability of a company.



### Determining No. of Clusters with Dendrogram

If you want to create flat clusters we can analyze the above dendrogram to determine no. of clusters. We first assume that the horizontal lines are extended on both sides, and as such, they would also cross the vertical lines. Now we have to identify the tallest vertical line that does not have any horizontal line crossing through it.

In the above dendrogram graph, such a vertical line is the blue line. We now draw a horizontal line across this vertical line as shown below. This horizontal line cuts the vertical line at two places, and this means the optimal number of clusters is 2.

Another way is to visually see which vertical line is showing the biggest jump. Since the vertical line denotes the distance or similarity between the two clusters, the big jump signifies the two clusters are not very similar. Again draw the horizontal line through this vertical line and the number of cuts it makes is optimal no. of clusters. Again in our example, it is the blue line, and the horizontal line cuts at two places so no. of clusters is 2.

(It should be noted however that these methods do not always guarantee the optimal number of clusters, it is just a guideline)

### Limitations of Agglomerative Clustering

49 stocks  $\rightarrow 49 \times 49$

1M data  $\rightarrow 1M \times 1M \rightarrow$  computationally intensive

space complexity  $\rightarrow O(N^2)$

Time complexity  $\rightarrow O(N^3)$

$(10^7)^3 \rightarrow$  hours

